



# Dictionary learning from phaseless measurements

Andreas Tillmann, Yonina Eldar, Julien Mairal

## ► To cite this version:

Andreas Tillmann, Yonina Eldar, Julien Mairal. Dictionary learning from phaseless measurements. IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP), Mar 2016, Shanghai, China. pp.4702-4706, 10.1109/ICASSP.2016.7472569 . hal-01387416

**HAL Id: hal-01387416**

**<https://inria.hal.science/hal-01387416>**

Submitted on 25 Oct 2016

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# DICTIONARY LEARNING FROM PHASELESS MEASUREMENTS

Andreas M. Tillmann<sup>\*</sup>

Yonina C. Eldar<sup>†</sup>

Julien Mairal<sup>°</sup>

<sup>\*</sup> TU Darmstadt, Research Group Optimization, Dolivostr. 15, 64293 Darmstadt, Germany

<sup>†</sup> Department of Electrical Engineering, Technion, Haifa 32000, Israel

<sup>°</sup> Inria, Lear Team, 655 Avenue de l'Europe, 38330 Montbonnot, France

## ABSTRACT

We propose a new algorithm to learn a dictionary along with sparse representations from signal measurements without phase. Specifically, we consider the task of reconstructing a two-dimensional image from squared-magnitude measurements of a complex-valued linear transformation of the original image. Several recent phase retrieval algorithms exploit underlying sparsity of the unknown signal in order to improve recovery performance. In this work, we consider sparse phase retrieval when the sparsifying dictionary is not known in advance, and we learn a dictionary such that each patch of the reconstructed image can be sparsely represented. Our numerical experiments demonstrate that our proposed scheme can obtain significantly better reconstructions for noisy phase retrieval problems than methods that cannot exploit such “hidden” sparsity.

**Index Terms**— Phase Retrieval, Dictionary Learning

## 1. INTRODUCTION

Phase retrieval has been an active research topic for decades [1, 2]. In mathematical terms, it can be formulated as

$$\text{find } \mathbf{x} \in \mathcal{X} \subseteq \mathbb{C}^N \quad \text{s.t.} \quad |f_i(\mathbf{x})|^2 = y_i \quad \forall i = 1, \dots, M, \quad (1)$$

where the functions  $f_i : \mathbb{C}^N \rightarrow \mathbb{C}$  are linear operators and the scalars  $y_i = |f_i(\hat{\mathbf{x}})|^2$  are nonlinear measurements of the original signal  $\hat{\mathbf{x}}$  in  $\mathcal{X}$ , obtained by removing the phase information. A traditional approach is to consider cases in which the solution of (1) is unique, up to global phase, and to devise algorithms to recover  $\hat{\mathbf{x}}$ . Usually, recovering  $\hat{\mathbf{x}}$  involves oversampling, i.e., taking  $M > N$  measurements. The most popular techniques for phase retrieval are based on alternating projections [3, 4, 5]. These methods generally require precise prior information on the signal (such as knowledge of the support set) and often converge to erroneous results. More recent approaches include semidefinite programming [6, 7, 8, 9, 10] and gradient-based methods such as Wirtinger Flow [11, 12].

In recent years, new techniques for (1) were developed when  $\hat{\mathbf{x}}$  is *sparse*, i.e., is composed of only few atoms from a known dictionary [7, 13, 12, 14], leading to algorithms with improved recovery performance. The main idea is akin to compressed sensing, where one works with fewer *linear* measurements than signal components [15, 16]. An important observation that boosted the applicability of sparse recovery techniques is that many classes of signals admit a sparse approximation in some basis or overcomplete dictionary [17, 18, 19]. While sometimes such dictionaries are known explicitly, better results can often be achieved by adapting the dictionary to the data [18]. Numerous algorithms have been developed for this task (e.g., [20, 21]) when the signal measurements are *linear*.

In this work, we propose a dictionary learning scheme for simultaneously solving the signal reconstruction and sparse representation

problems given *nonlinear*, *phaseless* and possibly *noisy* measurements. We demonstrate its ability to achieve significantly improved reconstructions (especially when the oversampling ratio is low and the noise level high) compared to Wirtinger Flow (WF) [11], which cannot exploit sparsity if the dictionary is unknown. Our algorithm, referred to as phase-retrieval dictionary-learning (PRDL), is based on alternating minimization where we iterate between best fitting the data and sparsely representing the recovered image. PRDL combines projected gradient descent like in Wirtinger Flow to update the image, iterative shrinkage to sparsify the image [22] and block-coordinate descent for the dictionary update [21].

## 2. PHASE-RETRIEVAL DICTIONARY-LEARNING

Our focus is on the application to image reconstruction. Therefore, we work in the 2D setting directly; however, all expressions and algorithms could easily be formulated for 1D signals, as in (1).

We wish to recover an image  $\hat{\mathbf{X}}$  in  $[0, 1]^{N_1 \times N_2}$  from noise-corrupted phaseless nonlinear measurements

$$\mathbf{Y} := |\mathcal{F}(\hat{\mathbf{X}})|^2 + \mathbf{N}, \quad (2)$$

where  $\mathcal{F} : \mathbb{C}^{N_1 \times N_2} \rightarrow \mathbb{C}^{M_1 \times M_2}$  is a linear operator,  $\mathbf{N}$  denotes noise, and the complex modulus and squares are taken component-wise. Signal sparsity is known to aid in phase retrieval, but a sparsifying transform is not always known a priori. This motivates learning a dictionary  $\mathbf{D}$  in  $\mathbb{R}^{s \times n}$  such that each  $s_1 \times s_2$  patch  $\hat{\mathbf{x}}^i$  of  $\hat{\mathbf{X}}$ , represented as a vector of size  $s = s_1 s_2$ , can be approximated by  $\hat{\mathbf{x}}^i \approx \mathbf{D} \mathbf{a}^i$  with a sparse vector  $\mathbf{a}^i$  in  $\mathbb{R}^n$ . Here,  $n$  is chosen a priori and the number of patches depends on whether the patches are overlapping or not. In general,  $\mathbf{D}$  is chosen such that  $n \geq s$ .

Before presenting our approach for tackling (2), we define the following notation. We consider the linear operator  $\mathcal{E} : \mathbb{C}^{N_1 \times N_2} \rightarrow \mathbb{C}^{s \times p}$  that extracts the  $p$  patches  $\hat{\mathbf{x}}^i$  (which may overlap or not) from an image  $\mathbf{X}$  and forms the matrix  $\mathcal{E}(\mathbf{X}) = (\mathbf{x}^1, \dots, \mathbf{x}^p)$ . Similarly, we define the linear operator  $\mathcal{R} : \mathbb{C}^{s \times p} \rightarrow \mathbb{C}^{N_1 \times N_2}$  that reverses this process, i.e., builds an image from a matrix containing vectorized patches as its columns. Thus, in particular, we have  $\mathcal{R}(\mathcal{E}(\mathbf{X})) = \mathbf{X}$ . Further, let  $\mathbf{A} := (\mathbf{a}^1, \dots, \mathbf{a}^p)$  in  $\mathbb{R}^{n \times p}$  be the matrix containing the patch representation coefficient vectors as columns. Then, our desired sparse-approximation relation  $\hat{\mathbf{x}}^i \approx \mathbf{D} \mathbf{a}^i$  can be expressed as  $\mathcal{E}(\hat{\mathbf{X}}) \approx \mathbf{D} \mathbf{A}$ .

With this notation in hand, we now introduce our method, called phase-retrieval dictionary-learning (PRDL), which aims at solving

$$\begin{aligned} \min_{\mathbf{X}, \mathbf{D}, \mathbf{A}} \quad & \frac{1}{4} \|\mathbf{Y} - |\mathcal{F}(\mathbf{X})|^2\|_F^2 + \frac{\mu}{2} \|\mathcal{E}(\mathbf{X}) - \mathbf{D} \mathbf{A}\|_F^2 + \lambda \sum_{i=1}^p \|\mathbf{a}^i\|_1 \\ \text{s.t.} \quad & \mathbf{X} \in [0, 1]^{N_1 \times N_2}, \quad \mathbf{D} \in \mathcal{D}. \end{aligned} \quad (3)$$

---

**Algorithm 1** Phase-Retrieval Dictionary-Learning Algorithm (PRDL)

---

**Input:** initial image estimate  $\mathbf{X}_{(0)} \in [0, 1]^{N_1 \times N_2}$ , initial dictionary  $\mathbf{D}_{(0)} \in \mathcal{D} \subset \mathbb{R}^{s \times n}$ , parameters  $\mu, \lambda > 0$ , iteration limits  $K_1, K_2$ ;

**Output:** Learned dictionary  $\mathbf{D} = \mathbf{D}_{(K)}$ , patch representations  $\mathbf{A} = \mathbf{A}_{(K)}$ , image reconstruction  $\mathbf{X} = \mathbf{X}_{(K)}$ ;

```
1: for  $\ell = 0, 1, 2, \dots, K_1 + K_2$  do
2:   choose step size  $\gamma_\ell^A$  as explained in Section 3.1 and update  $\mathbf{A}_{(\ell+1)} \leftarrow \mathcal{S}_{\lambda\gamma_\ell^A/\mu}(\mathbf{A}_\ell - \gamma_\ell^A \mathbf{D}_{(\ell)}^\top (\mathbf{D}_{(\ell)} \mathbf{A}_\ell - \mathcal{E}(\mathbf{X}_{(\ell)})))$ ;
3:   choose step size  $\gamma_\ell^X$  and update  $\mathbf{X}_{(\ell+1)} \leftarrow \mathcal{P}_\mathcal{X}(\mathbf{X}_{(\ell)} - \gamma_\ell^X (\Re(\mathcal{F}^*(\mathcal{F}(\mathbf{X}) \odot (|\mathcal{F}(\mathbf{X})|^2 - \mathbf{Y}))) + \mu \mathbf{R} \odot \mathcal{R}(\mathcal{E}(\mathbf{X}) - \mathbf{D}\mathbf{A})))$ ;
4:   if  $\ell < K_1$  then
5:     keep  $\mathbf{D}_{(\ell+1)} \leftarrow \mathbf{D}_{(\ell)}$ ;
6:   else
7:     set  $\mathbf{B} \leftarrow \mathcal{E}(\mathbf{X}_{(\ell)}) \mathbf{A}_{(\ell)}^\top$  and  $\mathbf{C} \leftarrow \mathbf{A}_{(\ell)} \mathbf{A}_{(\ell)}^\top$ ;
8:     for  $j = 1, \dots, n$  do
9:       if  $C_{jj} > 0$  then
10:        update  $j$ -th column:  $(\mathbf{D}_{(\ell+1)})_{\cdot j} \leftarrow \frac{1}{C_{jj}} (\mathbf{B}_{\cdot j} - \mathbf{D}_{(\ell)} \mathbf{C}_{\cdot j}) + (\mathbf{D}_{(\ell)})_{\cdot j}$ ;
11:       else
12:        reset  $j$ -th column: e.g.,  $(\mathbf{D}_{(\ell+1)})_{\cdot j} \leftarrow$  random  $\mathcal{N}(0, 1)$  vector (in  $\mathbb{R}^s$ );
13:       normalize  $(\mathbf{D}_{(\ell+1)})_{\cdot j} \leftarrow \frac{1}{\|(\mathbf{D}_{(\ell+1)})_{\cdot j}\|_2} (\mathbf{D}_{(\ell+1)})_{\cdot j}$ 
```

---

Here,  $\|\mathbf{X}\|_F$  denotes the Frobenius matrix-norm, which generalizes the Euclidean norm to matrices. The parameters  $\mu, \lambda > 0$  in the objective (3) provide a way to control the trade-off between the data fidelity term from the phase retrieval problem and the approximation sparsity of the image patches. To that effect, we use the  $\ell_1$ -norm, which is well-known to have a sparsity-inducing effect [23]. In order to avoid scaling disambiguities, we also restrict  $\mathbf{D}$  to be in the subset  $\mathcal{D}$  of matrices with unit- $\ell_2$ -norm columns, and assume  $n < p$  (otherwise, each patch is trivially representable by a 1-sparse vector  $\mathbf{a}^i$  by including  $\mathbf{x}^i / \|\mathbf{x}^i\|_2$  as a column of  $\mathbf{D}$ ).

### 3. ALGORITHMIC FRAMEWORK

Problem (3) is non-convex and difficult to solve, similarly to classical dictionary learning [24, 17, 20, 19]. Therefore, we adopt an algorithm that provides monotonic decrease of the objective. (In fact, convergence to a stationary point of (3) can be obtained under mild assumptions, using results from [25]; we cannot give details here due to space restrictions.)

The algorithmic framework we employ is that of *alternating minimization*: For each variable  $\mathbf{D}$ ,  $\mathbf{A}$  and  $\mathbf{X}$  in turn, we take one step towards solving (3) w.r.t. this variable alone, keeping the other two fixed. We summarize our method in Algorithm 1, where the superscript  $*$  denotes the adjoint operator (for a matrix  $\mathbf{Z}$ ,  $\mathbf{Z}^*$  is thus the conjugate transpose),  $\Re(\cdot)$  extracts the real part of a complex-valued argument, and  $\odot$  denotes the Hadamard (element-wise) product of two matrices. The algorithm also involves the classical soft-thresholding operator  $\mathcal{S}_\tau(\mathbf{Z}) := \max\{0, |\mathbf{Z}| - \tau\} \odot \text{sign}(\mathbf{Z})$  and the Euclidean projection  $\mathcal{P}_\mathcal{X}(\mathbf{Z}) := \max\{0, \min\{1, \mathbf{Z}\}\}$  onto  $\mathcal{X} := [0, 1]^{N_1 \times N_2}$ ; here, all operations are meant component-wise.

To avoid training the dictionary on potentially useless early estimates, the algorithm performs two phases—while the iteration counter  $\ell$  is smaller than  $K_1$ , the dictionary is not updated. We next explain the algorithmic steps in more detail.

#### 3.1. Updating the Patch Representation Vectors

Updating  $\mathbf{A}$  (i.e., (3) with  $\mathbf{D}$  and  $\mathbf{X}$  fixed at their current values) calls for decreasing the objective

$$\sum_{i=1}^p \left( \frac{1}{2} \|\mathbf{D}_{(\ell)} \mathbf{a}^i - \mathbf{x}_{(\ell)}^i\|_2^2 + \frac{\lambda}{\mu} \|\mathbf{a}^i\|_1 \right). \quad (4)$$

Since the objective here is separable, we can update all vectors  $\mathbf{a}^i$  in parallel independently. To do so, we choose to perform one step of the algorithm ISTA (see, e.g., [22]), which has a monotonic decrease property. ISTA is a gradient-based method which performs the following update for each  $i = 1, \dots, p$ :

$$\mathbf{a}_{(\ell+1)}^i = \mathcal{S}_{\lambda\gamma_\ell^A/\mu}(\mathbf{a}_{(\ell)}^i - \gamma_\ell^A \mathbf{D}_{(\ell)}^\top (\mathbf{D}_{(\ell)} \mathbf{a}_{(\ell)}^i - \mathbf{x}_{(\ell)}^i)).$$

This update involves a gradient descent step followed by soft-thresholding. The step size parameter  $\gamma_\ell^A$  can be chosen as  $1/L_A$ , where  $L_A$  is an upper bound on the Lipschitz constant of the gradient; here,  $L_A = \|\mathbf{D}_{(\ell)}\|_2^2$  would do, but a better strategy is to use a backtracking scheme to automatically update  $L_A$  [22].

Constructing  $\mathbf{A}_{(\ell+1)}$  from the  $\mathbf{a}_{(\ell+1)}^i$  as specified above is equivalent to Step 2 of Algorithm 1.

#### 3.2. Updating the Image Estimate

With  $\mathbf{D} = \mathbf{D}_{(\ell)}$  and  $\mathbf{A} = \mathbf{A}_{(\ell+1)}$  fixed, updating  $\mathbf{X}$  consists of decreasing the objective

$$\frac{1}{4} \|\mathbf{Y} - |\mathcal{F}(\mathbf{X})|^2\|_F^2 + \frac{\mu}{2} \|\mathcal{E}(\mathbf{X}) - \mathbf{D}\mathbf{A}\|_F^2 \quad (5)$$

$$\text{with } \mathbf{X} \in \mathcal{X} = [0, 1]^{N_1 \times N_2}.$$

This problem can be seen as a regularized version of the phase retrieval problem (with regularization parameter  $\mu$ ) that encourages the patches of  $\mathbf{X}$  to be close to the sparse approximation  $\mathbf{D}\mathbf{A}$  obtained during the previous iterate.

Our approach to decrease the value of the objective (5) is by a simple projected gradient descent step. In fact, for  $\mu = 0$ , this reduces to the so-called *Wirtinger Flow* method [11], but with necessary modifications to take into account the constraints on  $\mathbf{X}$  (real-valuedness and variable bounds  $[0, 1]$ ).

The gradient (matrix) of  $\varphi(\mathbf{X}) := \frac{1}{4} \|\mathbf{Y} - |\mathcal{F}(\mathbf{X})|^2\|_F^2$  with respect to  $\mathbf{X}$  can be computed as

$$\nabla \varphi(\mathbf{X}) = \Re(\mathcal{F}^*(\mathcal{F}(\mathbf{X}) \odot (|\mathcal{F}(\mathbf{X})|^2 - \mathbf{Y}))),$$

by using the chain rule (we omit a more detailed derivation due to space limitations). The gradient of  $\psi(\mathbf{X}) := \frac{\mu}{2} \|\mathcal{E}(\mathbf{X}) - \mathbf{D}\mathbf{A}\|_F^2$  is given by

$$\nabla \psi(\mathbf{X}) = \mu \mathcal{E}^*(\mathcal{E}(\mathbf{X}) - \mathbf{D}\mathbf{A}) = \mu \mathbf{R} \odot \mathcal{R}(\mathcal{E}(\mathbf{X}) - \mathbf{D}\mathbf{A}),$$

where  $\mathbf{R} \in \mathcal{X}$  has entries  $r_{ij}$  equal to the number of patches the respective pixel  $x_{ij}$  is contained in. Note that if the whole image is divided into a complete set of non-overlapping patches,  $\mathbf{R}$  will just be the all-ones matrix; otherwise, the element-wise multiplication with  $\mathbf{R}$  undoes the averaging of pixel values performed by  $\mathcal{R}$  when assembling an image from overlapping patches.

Finally, the gradient w.r.t.  $\mathbf{X}$  of the objective in (5) is  $\nabla\varphi(\mathbf{X}) + \nabla\psi(\mathbf{X}) \in \mathbb{R}^{N_1 \times N_2}$ , and the update in Step 3 of Algorithm 1 is indeed shown to be a projected gradient descent step. Typically, a backtracking strategy may be used for choosing the step size  $\gamma_\ell^X$ , as detailed in the section devoted to numerical experiments.

### 3.3. Updating the Dictionary

To update the dictionary, i.e., to approximately solve (3) w.r.t.  $\mathbf{D}$  alone, keeping  $\mathbf{X}$  and  $\mathbf{A}$  fixed at their current values, we employ one pass of a block-coordinate descent (BCD) algorithm on the columns of the dictionary [21]. The objective to decrease may be written as

$$\sum_{i=1}^p \|\mathbf{D}\mathbf{a}_{(\ell+1)}^i - \mathbf{x}_{(\ell+1)}^i\|_2^2 \quad \text{s.t.} \quad \mathbf{D} \in \mathcal{D}, \quad (6)$$

and the update rule given by Steps 4–13 corresponds exactly to one iteration of [19, Algorithm 11] applied to (6), which ensures the monotonic decrease of the objective function.

## 4. NUMERICAL EXPERIMENTS

We first describe some aspects of the concrete models and our implementation and then proceed to discuss our experiments.

### 4.1. Experimental Setup

We consider several linear operators  $\mathcal{F}$  corresponding to different types of measurements. We denote by  $\mathbf{F}$  the (normalized) 2D-Fourier transform, and introduce two complex Gaussian matrices  $\mathbf{G} \in \mathbb{C}^{M_1 \times N_1}$ ,  $\mathbf{H} \in \mathbb{C}^{M_2 \times N_2}$ , whose entries are i.i.d. samples from the distribution  $\mathcal{N}(0, \mathbf{I}/2) + i\mathcal{N}(0, \mathbf{I}/2)$ . Then, we experiment with the operators  $\mathcal{F}(\mathbf{X}) = \mathbf{G}\mathbf{X}$ ,  $\mathcal{F}(\mathbf{X}) = \mathbf{G}\mathbf{X}\mathbf{G}^*$ ,  $\mathcal{F}(\mathbf{X}) = \mathbf{G}\mathbf{X}\mathbf{H}^*$ , and the coded diffraction model

$$\mathcal{F}(\mathbf{X}) = \begin{pmatrix} \mathbf{F}(\overline{\mathbf{M}}_1 \odot \mathbf{X}) \\ \vdots \\ \mathbf{F}(\overline{\mathbf{M}}_m \odot \mathbf{X}) \end{pmatrix}, \quad \mathcal{F}^*(\mathbf{Z}) = \sum_{j=1}^m (\mathbf{M}_j \odot \mathbf{F}^*(\mathbf{Z})), \quad (7)$$

where the  $\mathbf{M}_j$ 's are admissible coded diffraction patterns (CDPs), see for instance [11, Sect. 4.1]; in our experiments we used ternary CDPs, such that each  $\mathbf{M}_j$  is in  $\{0, \pm 1\}^{N_1 \times N_2}$ . Note that the operator  $\mathbf{F}$  is implemented using fast Fourier transforms.

To reconstruct  $\hat{\mathbf{X}}$ , we choose an oversampling setting where  $M_1 = 4N_1$ ,  $M_2 = 4N_2$  and/or  $m = 2$ , respectively. Moreover, we corrupt our measurements with additive white Gaussian noise  $\mathbf{N}$  such that  $\text{SNR}(\mathbf{Y}, |\mathcal{F}(\hat{\mathbf{X}})|^2 + \mathbf{N}) = 10$  dB for the Gaussian-type, and 20 dB for CDP measurements, respectively. Note that this settings yields a quite heavy noise level for the Gaussian cases, and a relatively low oversampling ratio for the CDPs.

### 4.2. Implementation Details

We choose to initialize our algorithm with a simple random image  $\mathbf{X}_{(0)}$  in  $\mathcal{X}$  to demonstrate the robustness of our approach with respect to its initialization. Nevertheless, other choices are possible.

For instance, one may also initialize  $\mathbf{X}_{(0)}$  with a power-method scheme similar to that proposed in [11], modified to account for the required real-valuedness and box-constraints. The dictionary is initialized as  $\mathbf{D}_{(0)} = (\mathbf{I}, \mathbf{F}_D)$  in  $\mathbb{R}^{s \times 2s}$ , where  $\mathbf{F}_D$  is the 2D discrete cosine transform (see, e.g., [18]).

To update  $\mathbf{A}$ , we use the ISTA implementation from the SPAMS package [21] with its integrated backtracking line search (for  $L_A$ ).<sup>1</sup> Regarding the step sizes  $\gamma_\ell^X$  for the update of  $\mathbf{X}$  (Step 3 of Alg. 1), we adopt the following simple strategy: Whenever the gradient step leads to a reduction in the objective function value, we accept it; otherwise, we recompute the step with  $\gamma_\ell^X$  halved until a reduction was achieved (no more than 100 trials). Regardless of whether  $\gamma_\ell^X$  was reduced or not, we reset its value to  $1.68\gamma_\ell^X$  for the next round; the initial step size is 1.

Finally, we consider non-overlapping  $8 \times 8$  patches (though our code can also handle overlap). We run PRDL (Algorithm 1) with  $K_1 = 25$  and  $K_2 = 50$ ; the regularization/penalty parameter values can be read from Table 1 (there,  $m_Y$  is the number of elements of  $\mathbf{Y}$ ). We remark that these parameter values were empirically found to work well for the instances considered here but are not the outcome of rigorous parameter benchmarking, which is a subject of future investigations.

### 4.3. Computational Experiments

We test our method on a collection of typical (grayscale) test images used in the literature, namely cameraman, house and peppers of size  $256 \times 256$ , and lena, barbara, boat, fingerprint and mandrill of size  $512 \times 512$ . All experiments were carried out on a Linux 64-bit quad-core machine (2.8 GHz, 8 GB RAM) running Matlab R2015a.

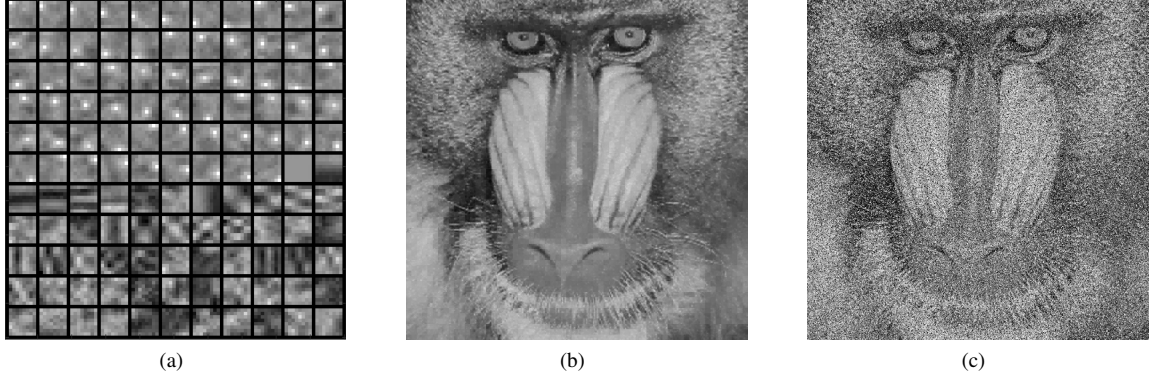
We evaluate our approach with the following question in mind: *Can we improve upon the quality of reconstruction compared to standard phase retrieval algorithms?* Standard methods cannot exploit sparsity if the underlying basis or dictionary is unknown; as we will see, the introduced (patch-)sparsity indeed allows for better recovery results (at least in the oversampling and noise regimes considered here).

To evaluate the achievable sparsity, we look at the average number of nonzeros in the columns of  $\mathbf{A}$  after running our algorithm. Generally, smaller values indicate an improved suitability of the learned dictionary for sparse patch coding (high values often occur if the regularization parameter  $\lambda$  is too small and the dictionary is learning the noise, which is something we would like to avoid). To assess the quality of the image reconstruction, we consider two standard measures, namely the peak signal-to-noise ratio (PSNR) of the reconstructed  $\mathbf{X}$  as well as its structural similarity index (SSIM) [26]. For PSNR, larger values are better, and SSIM-values closer to 1 (always ranging between 0 and 1) indicate better visual quality.

Table 1 displays the CPU times, PSNR- and SSIM-values and mean patch representation vector sparsity levels obtained for the various measurement types, averaged over the instance groups of the same size. The concrete example in Figure 1 shows the results from PRDL and plain Wirtinger Flow (WF; the real-valued,  $[0, 1]$ -box constrained variant, which corresponds to running Algorithm 1 with  $\mu = 0$  and omitting the updates of  $\mathbf{A}$  and  $\mathbf{D}$ ) for CDP measurements of the mandrill image. (In all tests, we let the Wirtinger Flow method run for the same number of iterations (75) and use the same starting points as for the PRDL runs.)

The PRDL method consistently provides better image reconstructions than WF, which clearly shows that our approach successfully introduces sparsity into the phase retrieval problem and exploits

<sup>1</sup><http://spams-devel.gforge.inria.fr/>.



**Fig. 1.** PRDL example: Image original is the  $512 \times 512$  “mandrill” picture, measurements are noisy CDPs (obtained using two ternary masks). (a) final dictionary (excerpt), (b) image reconstruction  $\mathcal{R}(\mathbf{DA})$  from sparsely coded patches, (c) reconstruction  $\mathbf{X}_{\text{WF}}$  after 75 WF iterations. Final PSNR 25.39 dB ( $\mathbf{X}_{\text{PRDL}}$ : 26.45,  $\mathbf{X}_{\text{WF}}$ : 13.42) and SSIM 0.7212 ( $\mathbf{X}_{\text{PRDL}}$ : 0.7601,  $\mathbf{X}_{\text{WF}}$ : 0.1562), average  $\|\mathbf{a}^i\|_0$  is 11.95.

$\mathcal{F}$ type	reconstruction	256 $\times$ 256 instances					512 $\times$ 512 instances				
		$(\mu, \lambda)/m_Y$	time	PSNR	SSIM	avg. $\ \mathbf{a}^i\ _0$	$(\mu, \lambda)/m_Y$	time	PSNR	SSIM	avg. $\ \mathbf{a}^i\ _0$
$G\hat{X}$	$\mathbf{X}_{\text{PRDL}}$	(0.5,0.068)	10.85	24.75	0.5370		(0.5,0.068)	52.93	21.18	0.5592	
	$\mathcal{R}(\mathbf{DA})$			24.66	0.6330	7.02			18.17	0.4964	2.12
	$\mathbf{X}_{\text{WF}}$		5.74	18.74	0.2812	–		29.72	18.80	0.3773	–
$G\hat{X}G^*$	$\mathbf{X}_{\text{PRDL}}$	(0.5,0.226)	42.55	22.42	0.4042		(0.5,0.068)	238.91	22.54	0.5275	
	$\mathcal{R}(\mathbf{DA})$			23.22	0.7270	6.72			23.82	0.7679	12.01
	$\mathbf{X}_{\text{WF}}$		32.53	22.41	0.4041	–		199.13	22.54	0.5274	–
$G\hat{X}H^*$	$\mathbf{X}_{\text{PRDL}}$	(0.5,0.226)	41.92	22.62	0.4114		(0.5,0.068)	242.71	22.53	0.5269	
	$\mathcal{R}(\mathbf{DA})$			23.24	0.7337	6.51			23.83	0.7688	11.99
	$\mathbf{X}_{\text{WF}}$		33.23	22.61	0.4111	–		202.40	22.53	0.5268	–
CDP (cf. (7))	$\mathbf{X}_{\text{PRDL}}$	(0.05,0.003)	7.51	27.20	0.7434		(0.05,0.003)	31.17	27.16	0.7764	
	$\mathcal{R}(\mathbf{DA})$			26.62	0.7675	7.79			26.18	0.7605	11.85
	$\mathbf{X}_{\text{WF}}$		2.29	13.03	0.1147	–		10.64	13.01	0.1539	–

**Table 1.** Test results for  $m_Y$  Gaussian-type and coded diffraction pattern (CDP) measurements. Reported are mean values (geometric mean for CPU times) per measurement type over the three  $256 \times 256$  and five  $512 \times 512$  instances w.r.t. the reconstructions from PRDL ( $\mathbf{X}_{\text{PRDL}}$  and  $\mathcal{R}(\mathbf{DA})$ ) with parameters  $(\mu, \lambda)$  and (real-valued,  $[0, 1]$ -constrained) Wirtinger Flow ( $\mathbf{X}_{\text{WF}}$ ), resp. (CPU times in seconds, PSNR in dB).

it in the solution process. As can be seen from Table 1, the obtained dictionaries allow for rather sparse representation vectors, with the effect of making better use of the information provided by the measurements, and also denoising the image along the way. The latter fact can be seen in the example (Fig. 1) and also inferred from the significantly higher PSNR and SSIM values for the estimates  $\mathbf{X}_{\text{PRDL}}$  and  $\mathcal{R}(\mathbf{DA})$  obtained from PRDL compared to the reconstruction  $\mathbf{X}_{\text{WF}}$  of the WF algorithm (which does not make use of sparsity).

## 5. DISCUSSION AND CONCLUSION

We believe that our experiments demonstrate that dictionary learning for phase retrieval with a patch-based sparsity is a promising direction, especially for cases where the original Wirtinger Flow approach fails. For the future, we are planning to conduct larger-scale experiments to answer some questions that we have left open.

For example, we expect further improvements by rigorously benchmarking the algorithmic parameters (in particular, to find good default settings). Several variants of our algorithm may also be developed—for instance, we successfully used  $\ell_0$ -constraints in-

stead of the  $\ell_1$ -penalty, by combining Orthogonal Matching Pursuit (OMP) [27] with our framework. Finally, to evaluate the quality of the learned dictionary, one might also ask how PRDL compares with the straightforward approach to first run (standard) phase retrieval and then learn dictionary and sparse patch representations from the result. Some preliminary experiments (not reported here) indicate that both approaches produce comparable results in the noise-free setting, while our numerical results demonstrate a denoising feature of our algorithm that the simple approach would obviously lack. Another interesting aspect to evaluate is by how much reconstruction quality and achievable sparsity degrade due to the loss of phase (or, more generally, measurement nonlinearity).

## Acknowledgements

The work of J. Mairal was funded by the French National Research Agency [Macaron project, ANR-14-CE23-0003-01]. The work of Y. Eldar was funded by the European Union’s Horizon 2020 research and innovation programme under grant agreement ERC-BNYQ, and by the Israel Science Foundation under Grant no. 335/14.

## 6. REFERENCES

- [1] S. Marchesini, “A Unified Evaluation of Iterative Projection Algorithms for Phase Retrieval,” *Review of Scientific Instruments*, vol. 78, pp. 011301–1–10, 2007.
- [2] Y. Shechtman, Y. C. Eldar, O. Cohen, H. N. Chapman, J. Miao, and M. Segev, “Phase Retrieval with Application to Optical Imaging,” *IEEE Signal Process. Mag.*, vol. 32, no. 3, pp. 87–109, 2015.
- [3] R. W. Gerchberg and W. O. Saxton, “A practical algorithm for the determination of the phase from image and diffraction plane pictures,” *Optik*, vol. 35, no. 2, pp. 237–246, 1972.
- [4] J. R. Fienup, “Phase Retrieval Algorithms: A Comparison,” *Appl. Opt.*, vol. 21, no. 15, pp. 2758–2769, 1982.
- [5] H. H. Bauschke, P. L. Combettes, and D. R. Luke, “Phase retrieval, error reduction algorithm, and Fienup variants: A view from convex optimization,” *J. Opt. Soc. Amer. A*, vol. 19, no. 7, pp. 1334–1345, 2002.
- [6] E. J. Candès, Y. C. Eldar, T. Strohmer, and V. Voroninski, “Phase retrieval via matrix completion,” *SIAM J. on Imaging Sciences*, vol. 6, no. 1, pp. 199–225, 2013.
- [7] Y. Shechtman, Y. C. Eldar, A. Szameit, and M. Segev, “Sparsity based sub-wavelength imaging with partially incoherent light via quadratic compressed sensing,” *Optics Express*, vol. 19, no. 16, pp. 14807–14822, 2011.
- [8] K. Jaganathan, S. Oymak, and B. Hassibi, “Recovery of sparse 1-D signals from the magnitudes of their Fourier transform,” arXiv:1206.1405[cs.IT], 2012.
- [9] H. Ohlsson, A. Y. Yang, R. Dong, and S. S. Sastry, “Compressive Phase Retrieval From Squared Output Measurements Via Semidefinite Programming,” arXiv:1111.6323[math.ST], 2012.
- [10] I. Waldspurger, A. d’Aspremont, and S. Mallat, “Phase recovery, MaxCut and complex semidefinite programming,” *Math. Prog. A*, vol. 149, no. 1, pp. 47–81, 2015.
- [11] E. J. Candès, X. Li, and M. Soltanolkotabi, “Phase Retrieval via Wirtinger Flow: Theory and Algorithms,” arXiv:1407.1065[cs.IT], 2014, to appear in IEEE Trans. Inf. Theory.
- [12] T. T. Cai, X. Li, and Z. Ma, “Optimal Rates of Convergence for Noisy Sparse Phase Retrieval via Thresholded Wirtinger Flow,” arXiv:1506.03382[math.ST], 2015.
- [13] Y. Shechtman, A. Beck, and Y. C. Eldar, “GESPAR: Efficient Phase Retrieval of Sparse Signals,” *IEEE Trans. Signal Process.*, vol. 62, no. 4, pp. 928–938, 2014.
- [14] M. L. Moravec, J. K. Romberg, and R. G. Baraniuk, “Compressive Phase Retrieval,” in *Proc. SPIE 6701: Wavelets XII*, D. Van De Ville, V. K. Goyal, and M. Papadakis, Eds. 2007, pp. 670120–1–11, International Society for Optical Engineering.
- [15] S. Foucart and H. Rauhut, *A Mathematical Introduction to Compressive Sensing*, Applied and Numerical Harmonic Analysis. Birkhäuser, 2013.
- [16] Y. C. Eldar and G. Kutyniok, Eds., *Compressed Sensing: Theory and Applications*, Cambridge University Press, 2012.
- [17] B. A. Olshausen and D. J. Field, “Emergence of Simple-Cell Receptive Field Properties by Learning a Sparse Code for Natural Images,” *Nature*, vol. 381, pp. 607–609, 1996.
- [18] M. Elad and M. Aharon, “Image Denoising Via Sparse and Redundant Representations Over Learned Dictionaries,” *IEEE Trans. Image Process.*, vol. 15, no. 12, pp. 3736–3745, 2006.
- [19] J. Mairal, F. Bach, and J. Ponce, “Sparse Modeling for Image and Vision Processing,” *Foundations and Trends in Computer Graphics and Vision*, vol. 8, no. 2-3, pp. 85–283, 2014.
- [20] M. Aharon, M. Elad, and A. Bruckstein, “K-SVD: An Algorithm for Designing Overcomplete Dictionaries for Sparse Representation,” *IEEE Trans. Signal Process.*, vol. 54, no. 11, pp. 4311–4322, 2006.
- [21] J. Mairal, F. Bach, J. Ponce, and G. Sapiro, “Online learning for matrix factorization and sparse coding,” *Journal of Machine Learning Research*, vol. 11, pp. 19–60, 2010.
- [22] A. Beck and M. Teboulle, “A Fast Iterative Shrinkage-Thresholding Algorithm for Linear Inverse Problems,” *SIAM J. Imaging Sci.*, vol. 2, no. 1, pp. 183–202, 2009.
- [23] S. S. Chen, D. L. Donoho, and M. A. Saunders, “Atomic Decomposition by Basis Pursuit,” *SIAM J. Sci. Comput.*, vol. 20, no. 1, pp. 33–61, 1998.
- [24] A. M. Tillmann, “On the Computational Intractability of Exact and Approximate Dictionary Learning,” *IEEE Signal Process. Lett.*, vol. 22, no. 1, pp. 45–49, 2015.
- [25] P. Tseng and S. Yun, “A coordinate gradient descent method for nonsmooth separable minimization,” *Math. Prog. B*, vol. 117, no. 1, pp. 387–423, 2009.
- [26] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, “Image Quality Assessment: From Error Visibility to Structural Similarity,” *IEEE Trans. Image Process.*, vol. 13, no. 4, pp. 600–612, 2004.
- [27] Y. C. Pati, R. Rezaifar, and P. S. Krishnaprasad, “Orthogonal Matching Pursuit: Recursive Function Approximation with Applications to Wavelet Decomposition,” in *Proc. 27th Annual Asilomar Conference on Signals, Systems and Computers*. 1993, vol. 1, pp. 40–44, IEEE Computer Society Press.